

EXHIBIT B

In the Matter of:
Berkeley J. Bates Seminar

Youtube Video

68 Commercial Wharf • Boston, MA 02110
888.825.3376 - 617.399.0130
Global Coverage
court-reporting.com



Google Exhibit 1022 Google v. Singular

Berkeley J. Bates Seminar

Youtube Video

COMMONWEALTH OF MASSACHUSETTS

YOUTUBE VIDEO

ON

PRACTICAL APPROXIMATE COMPUTING

AT

BERKELEY UNIVERSITY OF CALIFORNIA

SPEAKER: JOSEPH BATES

DATE: MARCH, 2016

MARY INDOMENICO, ACT, CET

official Court Transcriber

1 MR. BATES: I'm a computer scientist with a --
2 I don't know -- sore throat or something, so I hope you
3 can hear me.

4 I am currently at a start-up, small company.
5 Used to be at Carnegie Mellon AI faculty, and MIT, in
6 certain roles. And I want to talk about this
7 approximate computing idea, and what I think is a
8 practical simple approach to it.

9 So, I about ten years ago belatedly realized
10 what Carver Mead had been saying for a long time, which
11 is that wires -- you know, junctions can add currents,
12 and transistors can do exp and log. These are useful
13 operations in a lot of tasks. And a modern chip -- you
14 know, the chips in our cell phones, right, have
15 whatever billion transistors. They run whatever
16 gigahertz. So, in this simple sense of operation -- I
17 mean they're running digitally, but you know, they have
18 the capacity to do something like ten to the eighteenth
19 ops per second. That's probably more than at least my
20 brain is doing, best I understand neuroscience at least
21 at a -- you know, a certain level of abstraction.

22 But as a programmer, we only get essentially
23 zero of it. You know, its 0.0000, some more zeros
24 percent. One out of a billion, maybe; ten out of a

1 billion. And I'm an AI guy. I'm getting to be an
2 older AI guy. And I'd like to see my childhood dreams
3 come true before, you know, before its over for me.

4 So, one theory about what's made AI successful
5 to the extent it has been or it hasn't been, has been
6 the material scientists and physicists; no Einstein of
7 AI so far. Forgive me to any of my AI colleagues who
8 may be here. And so if that's true -- and deep
9 learnings, of course an example, right. I mean known
10 that in the 60s, known that's in the 80s, known that's
11 in the rest -- five years especially.

12 And so I, at least, think we or I desperately
13 need to see that compute in the world. I don't want it
14 going to make sure windows runs in the next generation
15 of laptops.

16 So, about a decade ago, I began to wonder
17 about this. And Carver said, you know, there's a way
18 you can do this. Look, just become a EE, not a
19 computer scientist, become a EE, but be a CMOS EE, and
20 make sure it's analoged and do it in sub threshold.

21 Thank you. Thanks a lot.

22 That turns out to be hard -- at least I found
23 it to be hard. I tried to do it. In fact, Jacob
24 helped me try to do it. But it seemed like there were

1 hundreds of people in the world skilled at this, and I
2 wanted there to be hundreds of thousands of people.
3 So, I asked myself the question: Could I take his
4 ideas and try to make a quote "normal" computer? And
5 si this is just going to tell you about a normal
6 computer built -- sort of inspired by these ideas.

7 So, suppose you look at machines that did
8 arithmetic that was pretty close. And the reason for
9 this is that, you know, multiply and divide circuits
10 are very large in silicon. So, I said look, you know,
11 I'm going to make life easier for the hardware people
12 and possibly somewhere between more difficult and
13 impossible for the software people. What would happen?
14 And so the spec I wrote down was not the I EEE 754
15 floating point spec, you know, like that. It was that.
16 That's the spec. People ask me: "Well what's the
17 distribution of the errors?" Don't think about it.
18 Don't program to it. Don't try to figure it out.
19 Okay? And the reason is going to be, you say this, the
20 hardware people get a lot of flexibility to do crazy
21 stuff to make tiny circuits. And you want tiny
22 circuits if you're in certain fields of science, 'cause
23 you want a lot of compute.

24 Now, the questions would be, you know: what

1 happens to the hardware? what happens to the software?
2 Is it impossible to program with a machine that
3 operates like that?

4 So, that's what I want to tell you about --
5 some of the results that I was able to find.

6 So, again, a lonely programmer was able -- me,
7 was able to take a floating point unit -- traditional
8 digital floating point unit and shrink it about 100X.
9 Now that's potentially very good -- potentially very
10 good. And we'll look at the issues. This was using
11 conventional digital silicon. And again, Gert helped
12 me back in the early days, and Jacob helped me. But I
13 couldn't figure out how to do it in analog.

14 So, I went -- you know, what people usually
15 do, you go back to digital silicon. But it's nice
16 'cause it's easy to fabricate. There's you know, and
17 hundreds of billions of dollars went into the fabs to
18 make these things. It's fast 'cause its digital. It's
19 deterministic, which programmers like because debugging
20 this, you know, for most people an important part of
21 their job function. And powering costs generally --
22 roughly scale with area. So, 100X smaller could be --
23 could be good on power and potentially cost. But you
24 can't have anything around that arithmetic unit. You

1 can't have a GPU's worth of control logic, because then
2 its pointless to make the arithmetic unit be little.

3 So, I looked for the simplest possible
4 hardware design I could find, and I took it as another
5 -- related go as you know, don't fight physics.
6 Distance is energy. And in the mathematized machines
7 that we all grew up with, there is no distance. Just
8 access it in memory, it's there, grab it. Okay?

9 And some people worry about cash, but most of
10 us don't try to figure out how the cashes work.

11 So, I wanted to expose that up to the
12 programmers because its really important if you want to
13 get the compute power that -- that the silicon
14 inherently has. You have to kind of face some reality.

15 So, the architecture I took was just the old
16 80s mesh connected computers where you have a -- and
17 this model looks like a SIMD model; it's -- the truth
18 is, the machine does MIMD, you know, every core can do
19 different things. But we're going to talk about it as
20 if its every core doing the same thing.

21 And forgive me, one thing I do -- I've been in
22 this thing for so long that -- like you know, I know
23 what a SIMD -- massively connected SIMD machine is and
24 how it works and all, but not everybody does. So,

1 through your training iterations, your epochs which are
2 very slow, the error falls down something like this.
3 But the interesting thing, of course, is that this
4 scale is epochs, not time or energy. And if you change
5 it to energy or time, the curve really looks like that.
6 And for the people that are doing -- spending a month,
7 they come in in the morning with an idea about a
8 network they want to -- us to trick. They start it, it
9 takes three weeks before they've got their training
10 results, and obviously they've forgotten largely what
11 they did and don't care about it anymore anyway -- it's
12 way too long. So, we need a lot more compute -- at
13 least in this deep learning field.

14 And the same thing works going forward. So,
15 you know, in an embedded system, a cell phone, or an
16 ear ring, or car, you can run these deep networks
17 efficiently.

18 Okay. So, briefly the hardware and then the
19 future and then I'll be done.

20 So, we built these chips, and they work.
21 Here's one -- I brought it actually -- works. They
22 have 2,000 cores in them. They're about a half a
23 centimeter on a side. That's very small -- relatively
24 small for a chip. A GPU is maybe twenty-five times

1 more area than this chip. This is done in a 40
2 nanometer silicon, which is a few generations old now.
3 You know, modern laptops will have 14 nanometer; that's
4 eight times denser silicon. These things run at about
5 a couple hundred megahertz. They get about 200
6 gigaflops per watt peak, which is a few times better
7 than the most recent GPUs that have come out. But
8 they're built using modern processes and this is built
9 using an older one.

10 We're putting sixteen of these chips on a
11 board. We're hooking it to a Lynx system, or an ARM.
12 We're building five of these for DARPA, so that's
13 170,000 core system. It produces about 68 teraflops
14 peak. More interesting is it has about 68 terabytes
15 per second memory bandwidth, which is a lot higher than
16 the corresponding GPU system. Because they're going
17 out currently over a bus to a separate memory. And
18 any time you -- you know, it's the opposite of
19 (indiscernible) local computing, and it kills -- it
20 kills us.

21 If you want to get the real power out of
22 silicon, I don't think you can use that kind of
23 architecture.

24 This -- this system's going on the net

1 probably this summer for people to explore. There are
2 some folks who have first dibs on it, but DARPA's being
3 very generous and saying anybody that has an
4 interesting project is welcome to try and use it.

5 And we also have these little embedded systems
6 that are available for certain groups.

7 This stuff -- the commercial people out here -
8 - this stuff is in patents.

9 So, there's a, you know, image of the chip.
10 My company -- another company in Boston -- Cadence, did
11 our physical design. MOSIS did -- and PW Run and
12 Global Foundries did the actual silicon.

13 Here's a modern GPU. That's a current high-
14 end GPU, desktop GPU done in 28 nanometer. Here's this
15 chip we built: 3,000 cores -- 2,000 cores. That one
16 has -- runs faster, but it's a hundred times the power
17 driving that machine.

18 And so the final slide here is we need to
19 continue building up our tools and libraries. I'm
20 working on a deep learning library to make it easier
21 for people in that field to try things. Building up
22 the community of users starting with the ones I showed
23 you in the prior slide. People doing real applications
24 this year I hope.

C E R T I F I C A T I O N

I, MARY INDOMENICO, AN APPROVED COURT
TRANSCRIBER, DO HEREBY CERTIFY THAT THE FOREGOING IS A
TRUE AND ACCURATE TRANSCRIPT FROM THE AUDIO RECORDING
PROVIDED TO ME BY THE OFFICE OF WOLF, GREENFIELD &
SACKS, P.C., IN THE PROCEEDINGS IN THE ABOVE ENTITLED
MATTER.

I, MARY INDOMENICO, FURTHER CERTIFY THAT THE
FOREGOING IS IN COMPLIANCE WITH THE ADMINISTRATIVE
OFFICE OF THE TRIAL COURT DIRECTIVE ON TRANSCRIPT
FORMAT.

I, MARY INDOMENICO, FURTHER CERTIFY THAT I
NEITHER AM COUNSEL FOR, RELATED TO, NOR EMPLOYED BY ANY
OF THE PARTIES TO THE ACTION IN WHICH THIS HEARING WAS
TAKEN, AND FURTHER THAT I AM NOT FINANCIALLY NOR
OTHERWISE INTERESTED IN THE OUTCOME OF THE ACTION.

Mary C. Indomenico

September 8, 2020

212 Vineland Avenue, East Longmeadow, MA 01028

413-746-1778

perfectinprint@aol.com

A handwritten signature in black ink, reading "Mary C. Indomenico", written over a horizontal line.